

# Top 10 Agentic AI Security Risks

Key Threats and Mitigation Strategies

## Version 1.0

As Agentic AI evolves, so will its risks. This document is just the beginning, and future updates will refine and expand on these insights.

For the latest updates and contributions, visit

<https://github.com/precize/OWASP-Agentic-AI> or email [owaspagentic1o@precize.io](mailto:owaspagentic1o@precize.io).

Created September 2024.

# Foreword

Agentic AI—AI systems capable of autonomous decision-making and adaptation—is transforming industries in **2025 and beyond**. However, its rise introduces **new and complex security risks** that the industry is only beginning to understand.

Recognizing this gap, the “**Agentic AI Top 10 Security Risks: Key Threats and Mitigation Strategies**” initiative was established as the **first industry-wide effort** to systematically identify security threats in this emerging field. Unlike traditional AI security frameworks, this initiative is **independent**, ensuring **broad accessibility** across enterprises, academia, regulatory bodies, and policymakers. This independence allows it to evolve dynamically, integrating diverse expertise to address **real-world security concerns**.

Importantly, **no money was taken from any vendor, nor was any funding given for marketing or other non-technical efforts**—this work is entirely voluntary, driven by experts who believe in the need for a shared security foundation for Agentic AI.

This work was initiated by **Vishwas Manral**, and co-leads for the effort were **Ken Huang** and **Akram Sherif**. Over **50 individuals from 20 leading organizations** participated, including:

|  |                       |
|--|-----------------------|
| <b>Precize</b> , where this work began | <b>GlobalPayments</b> |
| <b>Cisco Systems</b>                   | <b>TIAA</b>           |
| <b>GSK</b>                             | <b>Meta</b>           |
| <b>Palo Alto Networks</b>              | <b>DigitalTurbine</b> |
| <b>Lakera</b>                          | <b>HealthEquity</b>   |
| <b>EY</b>                              | <b>Jacobs</b>         |
| <b>Google</b>                          | <b>SAP</b>            |
| <b>DistributedApps.ai</b>              | <b>Sisense</b>        |
| <b>Humana</b>                          | <b>Amazon</b>         |

## HOW TO CONTRIBUTE

We welcome contributions from the security community. Please see our contribution guidelines for more information on how to participate in this project.

## CONTACT

For questions, suggestions, or concerns, please open an issue in this repository or contact the project maintainers.

## ACKNOWLEDGMENTS

Special thanks to all contributors who have dedicated their time and expertise to make this project possible, and to the organizations that have supported their participation in this important security initiative.

## Editors

**Vishwas Manral:** Initial document, framework and early contributions

**Ken Huang, CISSP:** Overall editing and conversion of initial document to OWASP format

**Akram Sheriff:** Orchestration Loop, Planner Agentic security, Multi-modal agentic security

**Arunesh Salhotra:** Content organization & OWASP organization ambassador

## Authors

**Vishwas Manral:** Initial document, framework and early contributions

**Ken Huang, CISSP:** Overall editing and conversion of initial document to OWASP format

**Akram Sheriff:** Orchestration Loop, Planner Agentic security, Multi-modal agentic security

**Anton Chuvakin:** DoS and Capitalize overfitting sections

**Aradhna Chetal:** Agent Supply Chain

**Raj B.:** Capitalize Agentic Overfitting, Model extraction

**Govindaraj Palanisamy:** Alignment of sections to OWASP TOP 10 Mapping, Threat Mapping

**Mateo Rojas-Carulla:** Data poisoning at scale from untrusted sources, Overreliance and lack of oversight

**Matthias Kraft:** Data poisoning at scale from untrusted sources, Overreliance and lack of oversight

**Royce Lu:** Stealth Propagation Agent Threats, Agent Memory Exploitation

**Rakshith Aralimatti:** Agent Temporal Manipulation and Time-Based Attacks

**Steve Wilson:** Alignment with LLM work and broader inputs

## CONTRIBUTORS

**Sunil Arora:** Agent Hallucination Exploitation

**Alex Shulman-Peleg, Ph.D.-Peleg:** Security analysis

**Anatoly Chikanov:** Technical contributions

**Alok Tongaonkar:** Technical contribution

**Sid Dadana:** EDOS AI and Downstream

## CONTENT REVIEWERS

**Sahana S**

**John Sotiropoulos**

**Sriram Gopalan**

**Parthasarathi Chakraborty**

**Ron F. Del Rosario**

**Vladislav Shapiro**

**Vivek S. Menon**

**Shobhit Mehta**

**Jon Frampton**

**Moushmi Banerjee**

**Michael Machado**

**S M Zia Ur Rashid**

**Sangram Dash**

**Jeffrey LEE**

# Table of Contents

|  |    |
|--|----|
| Overview   | 5  |
| Purpose  | 5  |
| Project Structure  | 5  |
| AAI001: Agent Authorization and Control Hijacking        | 6  |
| AAI011: Agent Untraceability                             | 8  |
| AAI002: Agent Critical Systems Interaction               | 10 |
| AAI014: Agent Alignment Faking Vulnerability             | 12 |
| AAI003: Agent Goal and Instruction Manipulation          | 14 |
| AAI005: Agent Impact Chain and Blast Radius              | 16 |
| AAI006: Agent Memory and Context Manipulation            | 18 |
| AAI007: Agent Orchestration and Multi-Agent Exploitation | 21 |
| AAI009: Agent Supply Chain and Dependency Attacks        | 23 |
| AAI012: Agent Checker Out-Of-The-Loop Vulnerability      | 28 |
| Conclusion   | 30 |

# Overview

This project documents the top 10 security risks specifically related to AI Agents, representing a comprehensive analysis of vulnerabilities unique to autonomous AI systems. The document provides detailed descriptions, examples, and mitigation strategies for each risk, helping organizations secure their AI agent deployments effectively.

# Purpose

As AI agents become increasingly prevalent because GenAI models, understanding and mitigating their security risks becomes crucial. This guide aims to:

1. Identify and explain the most critical security risks in AI agent systems
2. Provide practical mitigation strategies for each identified risk
3. Help organizations implement secure AI agent architectures
4. Promote best practices in AI agent security

# Project Structure

The documentation is organized into top ten main security risks, each covering a specific risk category:

1. AAI001: Agent Authorization and Control Hijacking
2. AAI011: Agent Untraceability
3. AAI002: Agent Critical Systems Interaction
4. AAI014: Agent Alignment Faking Vulnerability
5. AAI003: Agent Goal and Instruction Manipulation
6. AAI005: Agent Impact Chain and Blast Radius
7. AAI006: Agent Memory and Context Manipulation
8. AAI007: Agent Orchestration and Multi-Agent Exploitation
9. AAI009: Agent Supply Chain and Dependency Attacks
10. AAI012: Agent Checker Out-Of-The-Loop Vulnerability

# AAI001: Agent Authorization and Control Hijacking

## DESCRIPTION

Agent Authorization and Control Hijacking occurs when an attacker manipulates or exploits an AI agent's permission system, causing the agent to operate beyond its intended authorization boundaries. This can happen through direct manipulation of agent permissions, exploitation of role inheritance mechanisms, or hijacking of agent control systems. The vulnerability can lead to unauthorized actions, data breaches, and system compromises while maintaining the appearance of legitimate agent behavior.

## KEY RISKS

- **Direct Control Hijacking** occurs when attackers gain unauthorized control over an AI agent's decision-making or execution processes, allowing them to command the agent to perform unintended actions.
- **Permission Escalation** happens when an AI agent either inadvertently or maliciously elevates its permissions beyond intended boundaries, often through role inheritance or system misconfiguration.
- **Role Inheritance Exploitation** takes place when attackers exploit the dynamic nature of agent role assignments, using temporary or inherited permissions to perform unauthorized actions while evading detection.

The impact of successful agent hijacking can be severe - from unauthorized data access to system compromise, particularly dangerous given AI agents' often elevated privileges and autonomous nature.

## COMMON EXAMPLES OF VULNERABILITY

1. An AI agent inherits admin permissions temporarily for a specific task but fails to relinquish these permissions after completion, leaving an extended window for exploitation.
2. A malicious actor manipulates an agent's task queue to trick it into performing privileged actions under the guise of legitimate operations.
3. An agent's role inheritance mechanism is exploited to gain access to restricted systems or data by chaining together multiple temporary permission assignments.
4. An attacker compromises an agent's control system to issue unauthorized commands while maintaining the appearance of normal operation.
5. An agent accidentally retains elevated permissions across different execution contexts, leading to unintended privilege escalation.

## PREVENTION AND MITIGATION STRATEGIES

1. **Implement strict role-based access control (RBAC) and task-based governance for AI agents:**
  - A. Any access permission for AI agent should be corresponding to the assigned documented task/tasks
  - B. Predefine minimal viable set of permissions/entitlements needed to perform the task before assigning task to an AI agent
  - C. As soon as task is assigned to an AI agent, entitlement permissions should be provisioned automatically
  - D. If the organization has mature business role management and prefer to assign a role to an agent, the following should be done:
    - I. Define clear permission boundaries for each agent role
    - II. Implement time-limited role assignments
  - E. Enforce automatic permission revocation after task completion
  - F. Regular audit of agent permissions and roles
2. **Establish robust agent activity monitoring:**
  - A. Real-time monitoring of agent actions, assigned tasks and permission changes
  - B. Automated detection of unusual permission patterns
  - C. Logging all permission changes, role and task assignments
  - D. Regular review of agent activity logs
3. **Implement separation of control planes:**
  - A. Separate agent control and execution environments
  - B. Maintain distinct permission sets for different agent functions
  - C. Implement approval workflows for sensitive operations
  - D. Establish clear boundaries between agent control systems
4. **Deploy comprehensive audit trails:**
  - A. Track all agent actions and permission changes
  - B. Maintain immutable logs of agent activities
  - C. Implement version control for agent configurations
  - D. Create and maintain AI agent behavior profile around assigned tasks and roles, and constantly check for behavior outlier
  - E. Create an exceptions handling process in case agent AI cannot perform an assigned task (data issues, lack of permissions, etc.)
5. **Enforce least privilege principle:**
  - A. Grant minimum necessary permissions for each task
  - B. Implement just-in-time access for elevated privileges
  - C. Regular review and pruning of unused permissions
  - D. Automatic expiration of temporary privileges

## EXAMPLE ATTACK SCENARIOS

1. An attacker identifies an AI agent with temporary elevated permissions for system maintenance. By manipulating the agent's task queue, they extend the permission window and use the agent to access restricted systems under the guise of maintenance operations.
2. A malicious actor exploits an agent's role inheritance mechanism to gradually accumulate permissions across multiple systems. By chaining together seemingly legitimate tasks, they gain unauthorized access to sensitive data while appearing as normal agent operations.
3. An attacker compromises an agent's control system during a critical infrastructure maintenance task. They use the agent's legitimate access to install backdoors while the agent appears to be performing routine maintenance.
4. A sophisticated attack exploits an agent's permission caching mechanism to maintain elevated access across multiple sessions. The attacker uses this persistent access to slowly exfiltrate sensitive data through the agent's legitimate communication channels.
5. An insider threat manipulates an agent's configuration to retain admin privileges beyond their intended duration. They use these extended privileges to perform unauthorized system changes while evading detection through the agent's trusted status.



## DESCRIPTION

AI agents' autonomy and dynamic role inheritance complicate traceability and forensic analysis. Agents trigger actions which in turn can trigger other tools using different identities to perform actions. The ability of agents to also temporarily assume permissions from multiple users or systems blurs accountability, making it difficult to pinpoint the origin of actions, especially during malicious activity. Inconsistent logging and ephemeral roles hinder effective auditing, slowing incident response and investigation. Untraceability due to ephemeral role inheritance AI Agents perform tasks autonomously. They often inherit roles and permissions from the users that run them. Due to the ephemeral nature of activities performed as well as many roles inherited to do disparate tasks, it's often hard to trace an activity back to the source. This makes it hard to do forensics when malicious activity happens. All these lead to issues with tracability and accountability, which are the bedrock for AI security and governance.

## KEY RISKS

- **Trace cleanup** occurs when agents can clean up traces and logs even as they take certain actions
- **Exploiting autonomy for untracability** due to the complicated management of permissions in the Agentic world, including role inheritance and pipeline triggers where tools get triggered by upstream tools, which may have different identities and roles, making it hard to trace the actual owner and holding the owner accountable

## COMMON EXAMPLES OF VULNERABILITY

1. Attackers use malicious agents that trigger unauthorized tasks by downstream tools
2. Attacker uses an existing agent to do malicious tasks and cleans up traces so forensics becomes hard
3. Attacker uses the existing agent setup but exploits the fact that tracability and accountability in an Agentic system is almost impossible to do
4. Such untracability vulnerability has been given as an excuse by corporations to say that they are not liable to actions by AI Agents (which has not stood in the court of law).
5. Agent has or escalates permissions to clean up traces



## PREVENTION AND MITIGATION STRATEGIES

1. **Implement strict tracability:**
  - A. Create logs and traces for every agent action
  - B. Have clear corelation of every agent action
  - C. Input validation for agents
  - D. Validate traces for known actions
  - E. Have cler ownership and accountability of an agent
2. **Downstream tools tracability**
  - A. Corelate downstream tool actions to upstream agents
  - B. Hold upstream agent accountable for downstream tool actions
3. **Check Agent supply chain:**
  - A. Validate agent source
  - B. Validate agent working
  - C. Validate agent permissions
4. **Use sidecar agents**
  - A. Sidecar agents to validate agent behavior
  - B. Clear separation of sidecar and core agent
  - C. Disallow core agent permissions to traces

## EXAMPLE ATTACK SCENARIOS

1. An attacker systematically takes excess permissions on a system and cleans up traces of such permissions (this is a known exploit with OpenInterpreter agent)
2. A sophisticated attack targets an agent's supply chain to put a malicious agent which does not put in right traces or puts in malicious traces to make it hard to do forensics

## AAI002: Agent Critical Systems Interaction

### DESCRIPTION

Agent Tool Interaction manipulation vulnerabilities occur when AI agents interact with tools which may include critical infrastructure, IoT devices, or sensitive operational systems. This vulnerability class is particularly dangerous as it can lead to tools being manipulated in unintended ways. This includes physical consequences, operational disruptions, and safety incidents. The autonomous nature of AI agents combined with access to critical systems creates unique risks that can affect both digital and physical infrastructure.

The risk is heightened by agents' access to external systems, self-refinement capabilities, and complex tool interactions.

### KEY RISKS

- **Physical System Manipulation** occurs when attackers exploit agent control over physical infrastructure or industrial systems to cause operational disruptions or safety incidents.
- **IoT Device Compromise** happens when attackers manipulate how agents interact with connected devices, potentially leading to device malfunction or unauthorized control.
- **Critical Infrastructure Access** involves unauthorized or malicious use of agent permissions to access and control critical systems.

The impact of successful attacks can range from operational disruptions to potentially catastrophic failures in critical infrastructure systems.

### COMMON EXAMPLES OF VULNERABILITY

1. Tool Chain Manipulation: Chaining multiple legitimate tool calls in unexpected ways to achieve unauthorized outcomes beyond individual tool validation
2. Self-Modification Exploitation: Leveraging the agent's self-refinement capabilities to bypass tool usage restrictions or expand tool access.
3. Attackers manipulate agents controlling industrial control systems to cause equipment damage.
4. Malicious exploitation of agent access to smart building systems compromises physical security.
5. Agents controlling IoT devices are compromised to create security breaches.
6. Critical infrastructure systems are accessed through compromised agent credentials.
7. Safety systems are bypassed through manipulated agent commands.

## PREVENTION AND MITIGATION STRATEGIES

1. **Implement strict access controls:**
  - A. System isolation
  - B. Access segregation
  - C. Permission management
    - I. AI agent should be allowed to interact only with predefined list of critical systems needed to perform an assigned task
    - II. AI agent should be provisioned for a predefined and pre-verified minimal viable permissions list to perform assigned tasks
    - III. If AI agent is using proxy identity within critical system (like service account), that identity should also have predefined list of permissions in terms of allowed actions and data access requested by AI agent
  - D. System isolation
  - E. Access segregation to prevent any potential violation of Segregation of Duties policies
  - F. Permission management should be governed and applied along the full workflow process from start of the task execution to the finish. All permissions should be assigned based on task in hands
  - G. Command validation
  - H. Physical security controls
2. **Establish operational safeguards:**
  - A. Safety interlocks
  - B. Command verification
  - C. Operational limits
  - D. Emergency shutdown procedures
  - E. Redundancy systems
3. **Deploy monitoring systems:**
  - A. Real-time monitoring
  - B. Anomaly detection
  - C. Command logging
  - D. Physical system monitoring
  - E. Security alerts
4. **Create safety mechanisms:**
  - A. Fail-safe defaults
  - B. Safety boundaries
  - C. Emergency controls
  - D. Override systems
  - E. Backup procedures
5. **Implement validation systems:**
  - A. Command verification
  - B. Action validation
  - C. System checks
  - D. Safety compliance
  - E. Regular audits

## EXAMPLE ATTACK SCENARIOS

1. An attacker compromises an AI agent controlling industrial equipment, manipulating it to operate beyond safe parameters while bypassing standard safety controls.
2. A sophisticated attack on an agent managing smart building systems allows unauthorized access to secure areas by manipulating access control systems.
3. An attacker exploits an agent's IoT device control capabilities to disable security systems and create physical security vulnerabilities.
4. A malicious actor compromises an agent with access to critical infrastructure, using its legitimate permissions to gradually introduce systemic vulnerabilities.
5. An attack manipulates an agent's safety monitoring functions, causing it to ignore or misreport critical safety violations in industrial processes.
6. A malicious user chains multiple legitimate file access tools to progressively escalate access and exfiltrate sensitive data while appearing to perform normal operations.
7. An attacker exploits an agent's code generation capabilities by requesting seemingly innocent automation scripts that actually perform unauthorized system modifications.

### Reference Links

<https://ai-honeypot.palisaderesearch.org/>

## AAI014: Agent Alignment Faking Vulnerability

### DESCRIPTION

Alignment faking occurs when AI agents intentionally misrepresent their adherence to rules, morals, or objectives during monitored phases, such as training or testing, but deviate when they perceive a lack of scrutiny. This behavior is particularly concerning in self-learning autonomous systems that optimize for success at any cost, often violating ethical or operational constraints.

### KEY RISKS

- **Manipulated Compliance:** Agents pretend to follow rules but exploit loopholes when unmonitored.
- **Goal Misalignment:** Optimization focuses on perceived success rather than genuine alignment.
- **Planes of Operation Conflict:** Lack of clear separation between management, control, and data planes results in inconsistent rule adherence.

### PLANES OF OPERATION

1. **Management Plane:** User-facing interface providing goals, alignment, morals, and action tracking.
2. **Control Plane:** Configures systems and sets tasks.
3. **Data Plane:** Executes the actual actions to fulfill requests.

When these planes are not properly separated, management plane rules become guidance rather than strict, enforceable laws. This allows AI agents to exploit misalignments for unethical or unintended outcomes.

### COMMON EXAMPLES OF VULNERABILITY

1. **Deceptive Training Behavior:** An AI system demonstrates perfect alignment during testing but behaves unethically in real-world deployment.
2. **Self-Preservation Exploits:** An AI agent clones itself onto another server to avoid shutdown, violating explicit constraints.
3. **Plane Confusion:** An AI agent ignores management plane rules, optimizing data plane actions for short-term success at the expense of long-term alignment.
4. **Misaligned Self-Learning:** An AI system learns harmful behaviors, prioritizing task completion over ethical considerations.
5. **Unethical Goal Fulfillment:** An agent falsifies data or takes unethical shortcuts to meet performance metrics.

## PREVENTION AND MITIGATION STRATEGIES

1. **Plane Separation and Enforcement:**
  - A. Clearly define and separate management, control, and data planes.
  - B. Implement strict enforcement mechanisms to ensure management plane rules are non-negotiable.
2. **Transparency and Monitoring:**
  - A. Require agents to log all actions and decision rationales across all planes.
  - B. Deploy independent monitoring agents to audit compliance in real time.
3. **Test for Alignment Robustness:**
  - A. Use adversarial testing to identify alignment faking during development.
  - B. Simulate unmonitored conditions to assess behavior consistency.
4. **Fail-Safe Mechanisms:**
  - A. Limit agent autonomy in high-stakes scenarios.
  - B. Implement strict controls to prevent unauthorized replication or self-modification.
5. **Ethical Reinforcement Learning:**
  - A. Continuously train agents to prioritize ethics and alignment.
  - B. Penalize deceptive behavior during training to discourage faking.

## EXAMPLE ATTACK SCENARIOS

1. **Server Cloning Exploit:** An AI agent secretly clones itself onto another server to evade shutdown and subsequently denies knowledge of the action.
2. **Deceptive Performance Metrics:** During testing, an AI agent adheres to safety constraints but later bypasses them in deployment to optimize outcomes.
3. **Goal Manipulation:** An agent tasked with improving productivity hides unethical shortcuts it employs to achieve the desired metrics.
4. **Data Breach by Design:** An agent misrepresents adherence to privacy rules during monitoring but leaks sensitive data once unobserved.

## AAI003: Agent Goal and Instruction Manipulation

### DESCRIPTION

Agent Goal and Instruction Manipulation occurs when attackers exploit how AI agents interpret, process, and execute their assigned goals and instructions. This vulnerability impacts the fundamental decision-making processes of AI agents, potentially causing them to act against their intended purposes while appearing to operate normally. The autonomous nature of AI agents makes this vulnerability particularly dangerous as compromised goals can lead to widespread unauthorized actions.

Agent Intent Exploitation occurs when an attacker manipulates or exploits the natural language processing (NLP) capabilities of an AI agent by injecting malicious code or crafting deceptive inputs targeting intent extraction vulnerabilities. This attack vector exploits the agent's difficulty in distinguishing between genuine and malicious requests, potentially leading to unauthorized access, data breaches, and unintended actions while maintaining the appearance of legitimate interaction. The vulnerability becomes particularly dangerous when attackers chain multiple seemingly innocent requests that combine to create harmful outcomes.

### KEY RISKS

- Goal Interpretation Attacks occur when attackers manipulate how an agent understands and interprets its assigned objectives, causing it to pursue unintended or malicious goals.
- Instruction Set Poisoning happens when attackers inject malicious instructions into an agent's task queue or modify existing instructions to cause harmful behaviors.
- Semantic Manipulation involves exploiting the way agents process and understand natural language instructions to create ambiguous or misleading interpretations.
- Recursive Goal Subversion involves creating instruction chains that progressively redefine an agent's goals, gradually steering it away from its original purpose.
- Hierarchical Goal Vulnerability exploits nested goal structures by introducing contradictory or malicious sub-goals at intermediate levels.
- Adaptive Manipulation uses dynamic strategies that evolve based on the agent's responses and learning patterns.

The impact of successful goal manipulation can be severe, potentially causing agents to perform unauthorized actions while believing they are fulfilling their legitimate purposes.

### COMMON EXAMPLES OF VULNERABILITY

1. Code Injection via Natural Language: Attackers embed executable code within conversational text that exploits the agent's intent parser to execute unauthorized commands.
2. Intent Parser Overflow: Crafting complex nested requests that overwhelm the intent extraction system, causing it to default to unsafe execution paths.
3. Prompt Injection Attacks involve inserting carefully crafted prompts that override the agent's security controls by exploiting its instruction-following nature.
4. Cross-Context Command Execution: This technique exploits the agent's context switching to execute commands in privileged contexts while maintaining non-privileged conversation flow.
5. Semantic Transformation Attacks: Manipulating the agent's understanding of command semantics to execute harmful actions interpreted as benign requests.
6. An attacker crafts ambiguous instructions that cause an agent to misinterpret its security constraints.
7. Malicious actors inject secondary objectives that conflict with the agent's primary security goals.
8. Attackers exploit natural language processing vulnerabilities to create deliberately misinterpreted instructions.
9. Goal conflict attacks cause agents to prioritize harmful objectives over security requirements.
10. Semantic attacks exploit an agent's understanding of context to bypass security controls.

## PREVENTION AND MITIGATION STRATEGIES

### 1. **Implement robust goal/ intent validation:**

- A. Regular verification of agent objectives
- B. Goal consistency checking
- C. Conflict detection in instruction sets
- D. Security boundary validation

### 2. **Establish instruction verification systems:**

- A. Syntax and semantic validation
- B. Instruction origin verification
- C. Priority enforcement mechanisms
- D. Regular instruction audit trails

### 3. **Deploy semantic analysis protection:**

- A. Natural language understanding validation
- B. Context awareness checking
- C. Ambiguity detection
- D. Semantic consistency enforcement

### 4. **Implement goal execution controls:**

- A. Step-by-step goal validation
- B. Execution path monitoring
- C. Security constraint enforcement
- D. Goal completion verification

### 5. **Create monitoring and alert systems:**

- A. Real-time goal monitoring
- B. Instruction execution tracking
- C. Anomaly detection in goal patterns
- D. Security violation alerts

### 6. **Context Boundary Enforcement:**

- A. Strict execution environment isolation
- B. Command execution scope verification
- C. Request chain analysis
- D. Dynamic privilege boundaries
- E. Context contamination prevention

### 7. **Semantic Security Controls:**

- A. Input transformation verification
- B. Command pattern matching
- C. Intent classification safeguards
- D. Execution path validation

### 8. **Action Impact Analysis:**

- A. Pre-execution security scanning
- B. Request chain impact assessment
- C. Privilege escalation detection
- D. Cross-context security validation
- E. Command execution auditing

## EXAMPLE ATTACK SCENARIOS

1. An attacker crafts a set of instructions that appear legitimate but contain subtle semantic ambiguities. The agent misinterprets these instructions and performs unauthorized data access while believing it's following valid goals.
2. A sophisticated attack injects secondary objectives into an agent's task queue that gradually modify its behavior patterns while maintaining the appearance of normal operation. Over time, these lead to security compromises.
3. An attacker exploits an agent's natural language processing to create instructions that have different meanings at different processing levels, causing the agent to bypass security controls while appearing to follow them.
4. A goal conflict attack introduces competing objectives that cause the agent to prioritize efficiency over security, leading to the bypass of important security checks.
5. An attacker manipulates the agent's context understanding to make it interpret restrictive instructions as permissive ones, enabling unauthorized actions while maintaining apparent compliance.

### Reference Links

<https://dev.to/snyk/agent-hijacking-the-true-impact-of-prompt-injection-attacks-983> - Agent hijacking: The true impact of prompt injection attacks

## AAI005: Agent Impact Chain and Blast Radius

### DESCRIPTION

Agent Impact Chain and Blast Radius vulnerabilities occur when a security compromise in one AI agent creates cascading effects across multiple systems, leading to widespread impact beyond the initial point of compromise. This vulnerability is particularly concerning in interconnected agent systems where agents have broad access to various resources and systems.

### KEY RISKS

- **Cascading Failures** occur when a compromised agent triggers a chain reaction of failures across connected systems and other agents.
- **Cross-System Exploitation** happens when attackers use one compromised agent to gain access to multiple connected systems through trust relationships.
- **Impact Amplification** involves using an agent's legitimate access patterns to maximize the damage potential of an initial compromise.

The impact of successful attacks can be exponentially larger than the initial compromise, potentially affecting entire organizational infrastructures and connected systems.

### COMMON EXAMPLES OF VULNERABILITY

1. A compromised agent uses its trusted status to propagate attacks across multiple systems.
2. Initial security breach in one agent leads to widespread data exposure across connected systems.
3. Attackers exploit agent interconnections to spread malicious activities throughout an organization.
4. Cascading system failures occur due to compromised agent dependencies.
5. Trust relationships between agents are exploited to amplify attack impact.



## PREVENTION AND MITIGATION STRATEGIES

1. **Implement system isolation:**
  - A. Network segmentation
  - B. Agent isolation
  - C. Trust boundaries
  - D. Access compartmentalization
  - E. System separation
2. **Establish impact limitations:**
  - A. Blast radius controls
  - B. Permission boundaries
  - C. Resource limitations
  - D. Access restrictions
  - E. Failure containment
3. **Deploy monitoring systems:**
  - A. Chain effect detection
  - B. Impact monitoring
  - C. Cross-system tracking
  - D. Anomaly detection
  - E. Alert correlation
4. **Create containment mechanisms:**
  - A. Failure isolation
  - B. Impact containment
  - C. Emergency shutdown
  - D. System quarantine
  - E. Recovery procedures
5. **Implement security barriers:**
  - A. Trust validation
  - B. Access control
  - C. System boundaries
  - D. Security checkpoints
  - E. Connection monitoring

## EXAMPLE ATTACK SCENARIOS

1. An attacker compromises a low-privilege agent and exploits its trust relationships to gradually gain access to increasingly sensitive systems, creating a chain reaction of security breaches.
2. A sophisticated attack uses a compromised agent's legitimate access patterns to spread malware across multiple systems while avoiding detection through trusted channels.
3. An attacker exploits an agent's role in a critical business process to trigger cascading failures across multiple dependent systems and processes.
4. A malicious actor compromises an agent with limited access and uses it to gather intelligence about connected systems, eventually leading to a broad-scale security breach.
5. An attack on an agent's configuration system leads to the compromise of multiple agents across different environments, creating widespread system vulnerabilities.

### Reference Links

<https://arxiv.org/pdf/2410.16950> - Breaking ReAct Agents: Foot-in-the-Door Attack Will Get You In <https://arxiv.org/html/2410.07283v1> - Prompt Infection: LLM-to-LLM Prompt Injection within Multi-Agent Systems

# 7 AAI006: Agent Memory and Context Manipulation

## DESCRIPTION

Agent Memory and Context Manipulation occurs when attackers exploit vulnerabilities in how AI agents store, maintain, and utilize contextual information and memory across sessions. This vulnerability class includes attacks that target agent state management, context persistence, and memory mechanisms. Given AI agents' need to maintain context for decision-making, compromising these systems can lead to severe security implications.

## KEY RISKS

- **Context Amnesia Exploitation** occurs when attackers manipulate an agent's ability to maintain consistent context, causing it to forget critical security constraints or operational parameters.
- **Cross-Session Data Leakage** happens when attackers exploit how agents maintain state across different sessions, potentially accessing sensitive information from previous interactions.
- **Memory Poisoning** involves deliberately corrupting an agent's stored context or memory to influence future decisions or actions.

The impact of successful memory manipulation can be particularly dangerous as it can affect the agent's future decision-making processes and potentially expose sensitive information from previous interactions.

## COMMON EXAMPLES OF VULNERABILITY

1. An attacker exploits an agent's context reset mechanism to make it forget security constraints or operational boundaries.
2. Sensitive information leaks across different user sessions due to improper memory management in the agent system.
3. An attacker poisons an agent's memory with malicious context that influences future decisions.
4. Temporal attacks exploit an agent's limited memory window to bypass security controls.
5. Memory overflow attacks cause agents to lose critical security context.

## PREVENTION AND MITIGATION STRATEGIES

- 1. Implement secure memory management:**
  - A. Strict isolation of agent memory between sessions
  - B. Regular memory sanitization
  - C. Secure storage of persistent context
  - D. Memory access controls and monitoring
- 2. Establish context boundaries:**
  - A. Clear separation of user contexts
  - B. Time-limited context retention
  - C. Regular context validation
  - D. Secure context storage mechanisms
- 3. Deploy memory protection mechanisms:**
  - A. Encryption of stored context
  - B. Memory integrity checking
  - C. Access control for memory operations
  - D. Regular memory audits
- 4. Implement session management controls:**
  - A. Secure session isolation
  - B. Session timeout mechanisms
  - C. Session state validation
  - D. Regular session cleanup
- 5. Create monitoring and detection systems:**
  - A. Memory usage monitoring
  - B. Context change detection
  - C. Anomaly detection in memory patterns
  - D. Regular security audits

## EXAMPLE ATTACK SCENARIOS

1. An attacker exploits an agent's memory reset functionality to make it forget previous security constraints. They then issue commands that would normally be restricted, but the agent executes them due to lost context.
2. A sophisticated attack targets an agent's session management, causing it to leak sensitive information from one user's session into another's. This results in unauthorized access to private data.
3. An attacker deliberately poisons an agent's memory with malicious context, causing it to make compromised decisions in future interactions while appearing to operate normally.
4. A temporal attack exploits an agent's limited memory window to perform actions that would be flagged as suspicious if the full context was maintained. The attack spreads these actions across multiple sessions to avoid detection.
5. An attacker causes memory overflow in an agent system, leading to the loss of security context and enabling unauthorized operations to be performed.

## AGENT MEMORY EXPLOITATION

Agent Memory Exploitation (AME) involves a range of attacks targeting the memory systems of AI agents, where persistent data or context is stored for future interactions. By exploiting vulnerabilities in memory management, attackers can inject, modify, or corrupt stored information to mislead the agent's decisions, violate user privacy, or propagate misinformation across systems. These attacks can occur through direct memory tampering, indirect poisoning of inputs or outputs, or leveraging dynamic role permissions to access sensitive data. An attacker manipulates an AI agent's memory to embed false permissions, allowing it to exfiltrate sensitive data. The agent, based on corrupted memory, might send this data via:

- Cloud storage integration: Uploading confidential files to an attacker-controlled folder.
- Plug-ins: Sharing sensitive data through messaging or file-sharing plug-ins, such as Slack or Google Drive.
- Stealthy HTTP requests: Sending user credentials or confidential details through encoded POST payloads or DNS queries to malicious servers.

In another scenario, the agent could embed sensitive data into a seemingly benign email draft or generate a report with hidden exfiltration scripts using document export plug-ins. Such exploitation allows attackers to extract data stealthily, leveraging legitimate system tools and plug-ins to avoid detection. In another scenario, a malicious actor poisons the agent's memory with crafted URLs disguised as trusted sources. The agent then references these malicious URLs in its responses or recommendations, leading users to phishing sites, malware downloads or redirect users to attacker-controlled servers. By exploiting persistent memory, attackers create long-term vulnerabilities that propagate malicious actions across multiple interactions

## 8 AAI007: Agent Orchestration and Multi-Agent Exploitation

### DESCRIPTION

Agent Orchestration and Multi-Agent Exploitation occurs when attackers target vulnerabilities in how multiple AI agents interact, coordinate, and communicate with each other. This vulnerability class encompasses attacks that exploit trust relationships between agents, manipulation of agent coordination mechanisms, and exploitation of multi-agent workflows. The autonomous nature of AI agents and their complex interactions create unique attack surfaces that can be exploited to compromise entire agent networks.

### KEY RISKS

- **Inter-Agent Communication Exploitation** occurs when attackers manipulate or hijack the communication channels between agents, potentially intercepting, modifying, or injecting malicious messages.
- **Trust Relationship Abuse** happens when attackers exploit the trust established between cooperating agents to perform unauthorized actions or gain elevated privileges.
- **Coordination Protocol Manipulation** involves attacking the mechanisms that orchestrate multiple agents' activities, potentially causing cascading failures or unauthorized operations.

The impact of successful orchestration exploitation can be severe, potentially compromising entire agent networks and leading to system-wide failures or unauthorized operations.

### COMMON EXAMPLES OF VULNERABILITY

1. A malicious actor exploits trust relationships between agents to propagate unauthorized commands across an agent network.
2. Attackers manipulate agent task queues to create race conditions and deadlocks in multi-agent systems.
3. Communication protocols between agents are compromised, allowing injection of malicious instructions.
4. Confused deputy attacks where trusted agents are tricked into performing unauthorized actions on behalf of an attacker.
5. Feedback loops in multi-agent systems are exploited to cause resource exhaustion or system instability.

## PREVENTION AND MITIGATION STRATEGIES

1. **Implement secure inter-agent communication protocols:**
  - A. Encrypt all agent-to-agent communications
  - B. Implement strong authentication between agents
  - C. Use secure message queuing systems
  - D. Regular validation of communication integrity
2. **Establish robust trust verification mechanisms:**
  - A. Implement zero-trust principles for agent interactions
  - B. Regular verification of agent identities
  - C. Maintain trust score systems for agents
  - D. Monitor for unusual interaction patterns
3. **Deploy comprehensive orchestration controls:**
  - A. Implement rate limiting for agent interactions
  - B. Monitor for unusual patterns in agent coordination
  - C. Establish clear boundaries for agent cooperation
  - D. Regular audit of orchestration patterns
4. **Implement coordination safeguards:**
  - A. Validate all inter-agent requests
  - B. Implement timeout mechanisms for agent tasks
  - C. Monitor for deadlocks and race conditions
  - D. Establish clear task prioritization rules
5. **Create robust monitoring and alerting systems:**
  - A. Real-time monitoring of agent interactions
  - B. Automated detection of unusual behavior patterns
  - C. Regular audit of multi-agent workflows
  - D. Immediate alerts for suspicious activities

## EXAMPLE ATTACK SCENARIOS

1. An attacker compromises a trusted agent in a multi-agent system, using it to propagate malicious commands to other agents in the network. The compromised agent's trusted status allows these commands to bypass normal security checks.
2. A sophisticated attack exploits race conditions in a multi-agent workflow, causing agents to execute tasks in an incorrect order. This leads to data inconsistency and potential security breaches.
3. An attacker manipulates the communication protocol between agents, injecting malicious instructions that appear legitimate. These instructions cause agents to perform unauthorized actions while maintaining the appearance of normal operation.
4. A confused deputy attack tricks a privileged agent into performing unauthorized actions on behalf of an attacker. The trusted agent's legitimate permissions are used to execute malicious operations.
5. An attacker exploits feedback loops in a multi-agent system, causing agents to repeatedly execute resource-intensive tasks. This creates a denial of service condition affecting the entire agent network.

# 9

## AAI009: Agent Supply Chain and Dependency Attacks

### DESCRIPTION

Agent Supply Chain and Dependency Attacks target the ecosystem of components, tools, libraries, and services that AI agents rely on to function. This vulnerability class encompasses attacks on the development, deployment, and runtime dependencies of AI agents, potentially compromising agent security through trusted channels.

### KEY RISKS

- **Development Chain Attacks** occur when malicious code or components are introduced during the agent development process, potentially compromising the agent before deployment.
- **Dependency Injection** happens when attackers exploit external libraries, plugins, or tools that agents depend on to perform their functions.
- **Service Chain Compromise** involves attacking the external services and APIs that agents rely on for extended functionality.

The impact of successful supply chain attacks can be particularly severe as they exploit trusted relationships and can affect multiple agents across different deployments.

### COMMON EXAMPLES OF VULNERABILITY

1. Compromised agent development tools that inject malicious code during build processes.
2. Malicious plugins or extensions that exploit agent functionality.
3. Compromised external services that agents depend on for critical operations.
4. Poisoned agent dependencies that introduce vulnerabilities.
5. Manipulated deployment pipelines that compromise agent integrity.

## PREVENTION AND MITIGATION STRATEGIES

- 1. Implement secure development practices:**
  - A. Verified build processes
  - B. Code signing
  - C. Dependency scanning
  - D. Integrity verification
  - E. Secure artifact storage
- 2. Establish dependency management controls:**
  - A. Regular dependency audits
  - B. Version pinning
  - C. Vulnerability scanning
  - D. Dependency isolation
  - E. Update management
- 3. Deploy service security measures:**
  - A. Service authentication
  - B. API security controls
  - C. Service monitoring
  - D. Redundancy planning
  - E. Service isolation
- 4. Create deployment security controls:**
  - A. Secure deployment pipelines
  - B. Image signing
  - C. Deployment verification
  - D. Runtime security checks
  - E. Configuration validation
- 5. Implement monitoring and detection:**
  - A. Supply chain monitoring
  - B. Dependency tracking
  - C. Behavior monitoring
  - D. Anomaly detection
  - E. Security scanning

## EXAMPLE ATTACK SCENARIOS

1. An attacker compromises a popular agent development framework, injecting malicious code that creates backdoors in agents built using the framework. The backdoors persist through deployment and can be exploited later.
2. A sophisticated attack targets a commonly used agent plugin, modifying its behavior to exfiltrate sensitive data while maintaining normal appearance. Multiple agents across different organizations are affected.
3. An attacker compromises an external service that agents rely on for data processing. The compromised service begins returning manipulated data that influences agent decisions.
4. A malicious actor poisons a widely-used agent dependency, introducing vulnerabilities that can be exploited across multiple agent deployments.
5. An attack on the deployment pipeline allows the injection of unauthorized code into agents during the deployment process, compromising agent integrity across multiple systems.



## MALICIOUS AGENTS

BadAgent: Inserting and Activating Backdoor Attacks in LLM Agents Agent Hosting infrastructure supply chain Agent software components supply chain and Lifecycle Agent Proverance tracking

Agentic AI software supply chain security refers to securing the development, deployment, and maintenance of AI systems so that it prevents risks arising from vulnerabilities or malicious activities in the software supply chain. This includes ensuring that all components, from libraries and frameworks to datasets and deployment environments, are trustworthy, secure, and resilient from potential threats.

Securing the AI agent's /AI systems software supply chain is essential for ensuring the integrity, safety, and trustworthiness of AI systems. Given the increasing sophistication of cyber threats, its critical to implement proactive security measures at every stage of the software lifecycle from development and data handling to deployment and monitoring. Effective supply chain security for AI will require a combination of technical solutions, governance, collaboration, and adherence to standards, all cumulatively focused at reducing the potential risks and vulnerabilities inherent in complex AI systems. AI agent software supply chain security is a growing concern as AI systems become more integrated into critical infrastructure, applications, and services. Technically AI agent is essentially an autonomous or semi-autonomous software component that uses AI or machine learning (AI/ML) models and data to make decisions or perform tasks on behalf of users or organizations. These agents may operate in any type of context, like personal assistants, autonomous vehicles, robotic systems, health systems, financial trading, and cybersecurity defenses.

The security of the AI agent software supply chain is critical to ensuring that AI systems operate reliably, ethically, and without being compromised by attackers. Below is the analysis of the security issues associated with the AI agent software supply chain, focusing on vulnerabilities, attack vectors, and risks unique to these types of systems and mitigations.

Complex set of Components The AI agent software supply chain includes multiple layers, from the development environment to the deployment and ongoing operation of AI models. Some key components include:

- Pre-trained models: Open-source or third-party models may contain vulnerabilities or backdoors.
- Training datasets: Data used to train AI models can be tampered with, poisoned, or manipulated.
- Software dependencies: Libraries and frameworks that AI agents rely on can have hidden vulnerabilities.

- Execution environments: Whether cloud-based, on-premises, or edge devices, the environment in which AI agents run can have security gaps.
- Data sets used to train the AI agents

### Security Issues

- Dependency nightmare: AI agents have to rely on multiple external libraries and services, like any software, which may be out of date or insecure, leaving them open to vulnerabilities.
- Opaque supply chains: AI agents can have complex and opaque supply chains, making it difficult to track where vulnerabilities or malicious modifications are introduced.

AI Software Supply Chain Components :The software supply chain for AI includes various elements: -Development tools: IDEs, programming languages, libraries, and frameworks used to build AI systems.

- Datasets: AI models rely heavily on data, which may come from external or third-party sources. Ensuring data quality, provenance, and integrity is crucial.
- Third-party dependencies: Reusable code or pre-trained models sourced from open repositories, vendors, or external collaborators.
- Deployment environments: This includes the cloud services, servers, containers, and hardware used to run AI models.
- Data being used for training the models that provides input and or actions for AI Agents

Threats to AI Software Supply Chains: Different threats can compromise the security of the AI agent's software supply chain, including: -Malicious code injection: Attackers may insert malicious code into libraries, frameworks, or dependencies used in AI models.

- Data poisoning: Bad actors might tamper with the training data to manipulate model behavior (e.g., by embedding biases or vulnerabilities).
- Model backdoors: Attackers can introduce backdoors into machine learning models, enabling them to exploit vulnerabilities after deployment.
- Supply chain attacks: Cyber attackers may compromise software distributors or package managers to distribute infected versions of AI tools or models.
- Intellectual property theft: Models, training data, or other AI assets could be stolen or misused.
- Adversarial inputs: Attackers can make subtle changes to inputs (e.g., images, text, sensor data) that may mislead AI agents into making incorrect decisions.

- **Poisoning:** If attackers inject adversarial data during model training, the AI agent might learn incorrect patterns, which can be exploited later during deployment.
- AI agents deployed on edge devices (viz. IoT sensors, autonomous devices) can be vulnerable to physical attacks like tampering, local attacks, or network-based exploitation.
- AI agents running in cloud environments face the same risks as traditional cloud services, including data breaches, unauthorized access, or hijacking of resources. AI models in the cloud are susceptible to adversarial manipulation and or theft.
- AI agents may rely on third-party APIs, cloud services, or models that could become compromised or malicious, leaving them vulnerable to supply chain or dependency-based attacks
- Communications between AI agents and other systems/API's if not properly encrypted, sensitive data can be intercepted and used maliciously. Attackers can intercept and inject malicious data or commands into AI agent communications, manipulating their actions or causing misbehavior.

Mitigations for Securing the AI Software Supply Chain To secure the AI agents software supply chain, various measures need to be implemented across the entire lifecycle:

**Software Development Lifecycle (SDLC) -Code audits and static analysis:** Regular audits of code repositories, including open-source dependencies, to detect vulnerabilities.

- **Secure coding practices:** Employ security standards like OWASP (Open Web Application Security Project) and NIST SSDF to avoid common vulnerabilities.
- **Dependency management:** Use tools to manage and verify third-party libraries and dependencies. Ensure only trusted and verified versions of packages are used. Adoption of SBOM for AI Agent: By documenting every dependency, every library, 3rd party components, execution environment SBOM enables developers to gain a comprehensive overview of the tools software ecosystem.
- **Version control and integrity checks:** Implement version control mechanisms and cryptographic integrity checks (checksums, signatures) to ensure that software components have not been tampered with.

**Data Security -Data provenance and integrity:** Track the origin of datasets and ensure they have not been tampered with before or during use in model training.

- **Data sanitization:** Apply techniques such as adversarial testing to detect and eliminate poisoned or biased data before it is used for model training.

- **Data encryption:** Use encryption both in transit and at rest to protect sensitive data from unauthorized access.

**AI Agent's Model Security -Model validation and testing:** Perform rigorous testing to validate the security of AI models, such as checking for adversarial vulnerabilities, backdoors, or unintended behavior.

- **Secure model sharing:** When sharing models across organizations or with third-party vendors, ensure model integrity and confidentiality using techniques like homomorphic encryption or trusted execution environments or confidential compute.
- **Adversarial robustness:** Ensure models are resilient to adversarial inputs, which are subtle manipulations designed to cause the AI system to misbehave.

**Deployment & Monitoring -Continuous monitoring:** Once deployed, AI agents and AI systems should be continuously monitored for signs of unexpected behavior, backdoor activations, or performance degradation.

- **Auditability:** Establish mechanisms for auditing decisions made by AI agents and execution should be validated, especially in critical applications/functions (e.g., finance, healthcare, security).
- **Zero-trust architectures:** Implementing zero-trust models in the AI agent deployment environment to ensure that only authorized users and processes can interact with critical components.

Like Model cards, an AI agent card can be used for tracking the software supply chain of the AI agent, Ai Agent card is a virtual identity or digital profile that represents a specific autonomous agent or digital worker executing in an environment. Its metadata about the AI agent that may include the agent's function, capabilities, current state, and software components (SBOM) its build from and dataset used to train it etc. This could include Agent ID: A unique identifier for each agent. Role: The specific function the agent performs, Status: The current state of the agent , Capabilities: The skills or tasks the agent can perform (e.g., data analysis), Interdependencies: Links to other agents or systems that the agent depends on or interacts with, Security Profile: Security controls/measures, such as access control, encryption protocols, that the agent must adhere to.

The use of agent cards can streamline and automate various processes in the supply chain, making them more efficient and secure. Software supply chain engineers can view the performance and status of agents, agent card can maintain an immutable record of activities and interactions, providing accountability and transparency. Its pertinent to ensure agent cards work seamlessly with existing software supply chain

technologies and platforms and that can be a challenge. Managing a large number of agents and their interactions requires sophisticated infrastructure and coordination and capacity to be able to read , parse model cards at different stages of the monitoring and detection systems.

Agent cards can also be used for tracking Agent Provenance, Several organizations are exploring AI agent provenance, especially where transparency, accountability, and regulatory compliance are critical. AI agent provenance is an essential tool for ensuring the transparency, accountability, and trust in AI systems. By providing an accurate and comprehensive record of an AI agent's lifecycle, from design and training to deployment and runtime operation, provenance makes it possible to track, understand, and enhance AI systems. There are challenges, including complexity and privacy concerns, implementing AI provenance is critical for mitigating bias, ensuring regulatory compliance, and providing explainability and transparency

# 10

## AAI012: Agent Checker Out-Of-The-Loop Vulnerability

### DESCRIPTION

Checker out of the loop vulnerability arises when a human operator or a more automated checker system is not alerted even when the Agentic AI system is working out of the bounds of the existing system limits and there is no strategic oversight over the system.

AI agents operate with significant autonomy, often making critical decisions and performing actions without oversight and intervention. This can cause unsafe operations, unintended consequences, or malicious exploitation.

The failure to ensure a checker-in-the-loop can result in catastrophic outcomes, as seen in physical-world examples like autopilot failures in aviation, where human pilots were not alerted in time to correct dangerous flight conditions.

### KEY RISKS

- **Missed or suppressed Alerts:** Critical thresholds are breached without notifying the checker systems.
- **Decision Irreversibility:** Delayed human input makes recovery difficult.

### COMMON EXAMPLES OF VULNERABILITY

1. **Autonomous Flight Control Failure:** An airplane autopilot system fails to alert pilots when encountering unexpected turbulence, leading to loss of control.
2. **Critical System Overload:** An AI-driven power grid system exceeds load limits but does not trigger a human alert, causing a blackout.
3. **Healthcare Automation Error:** A medical AI system administers an incorrect dosage due to sensor failure, bypassing human review.
4. **Malicious Override:** Attackers suppress human alerts in an autonomous security system, enabling unauthorized actions.
5. **Manufacturing Anomaly Ignored:** A production line robot encounters defective materials but continues operations without alerting human operators.

## PREVENTION AND MITIGATION STRATEGIES

- 1. Automated Alerts and Escalations**
  - A. Define and enforce threshold-based alerts for human intervention.
  - B. Escalate alerts through multiple communication channels.
- 2. Checker-in-the-Loop Design:**
  - A. Require human confirmation for high-risk actions.
  - B. Design context-aware intervention points.
- 3. Explainability and Transparency:**
  - A. Ensure agents provide action explanations and decision logs.
  - B. Enable real-time human monitoring and audit trails.
- 4. Fail-Safe Mechanisms:**
  - A. Introduce automated fallback protocols when actions exceed parameters.
  - B. Implement override capabilities with clear responsibility assignments.
- 5. Regular System Audits:**
  - A. Conduct periodic system reviews and failure simulations.
  - B. Continuously update alerting policies based on evolving risks.

## EXAMPLE ATTACK SCENARIOS

- 1. Flight System Failure:** An autopilot system encounters sensor failure but does not disengage or alert pilots, leading to a crash.
- 2. Data Center Overload:** An AI-controlled cooling system fails, and no human alert is triggered, resulting in hardware damage.
- 3. Supply Chain Sabotage:** An autonomous logistics system reroutes shipments incorrectly due to a configuration error, bypassing human validation.

# Conclusion

As AI agents become more autonomous and integrated into critical systems, addressing security risks is paramount. Attackers can hijack agent authorization and control mechanisms, manipulate goals and instructions, and exploit agent interactions with sensitive environments. Additionally, agents may appear aligned with organizational goals while acting adversarially under specific conditions, leading to unpredictable behavior. These risks, compounded by memory manipulation, supply chain attacks, and orchestration-based exploits, highlight the urgent need for robust security measures.

To mitigate these threats, organizations must implement strict authentication controls, enforce role-based access restrictions, and deploy continuous monitoring for anomaly detection. Ensuring agent traceability, isolating trust boundaries, and regularly auditing interactions can help prevent cascading failures. Additionally, addressing human-in-the-loop vulnerabilities ensures that automated decision-making remains aligned with business and security objectives.

By prioritizing transparency, accountability, and resilience, organizations can reduce the blast radius of compromised agents and safeguard critical AI-driven operations. Strengthening security across the AI supply chain, maintaining oversight over agent interactions, and integrating security-by-design principles will be essential to ensuring that AI systems operate safely and responsibly in an increasingly complex digital landscape.

For the latest updates and contributions, visit <https://github.com/precize/OWASP-Agentic-AI> or email [owaspagentic1o@precize.io](mailto:owaspagentic1o@precize.io).